

Important Announcements

- assignment 1 is due this Friday (Jan 24)
- review session by TA on Monday
- notation : random variable \rightarrow upper case , X
outcome space / set
of a r.v. \rightarrow calligraphic
uppercase , \mathcal{X}
outcome of a r.v. \rightarrow lower cas , x

$$X \in \mathcal{X} = \{x_1, x_2, \dots, x_n\}$$

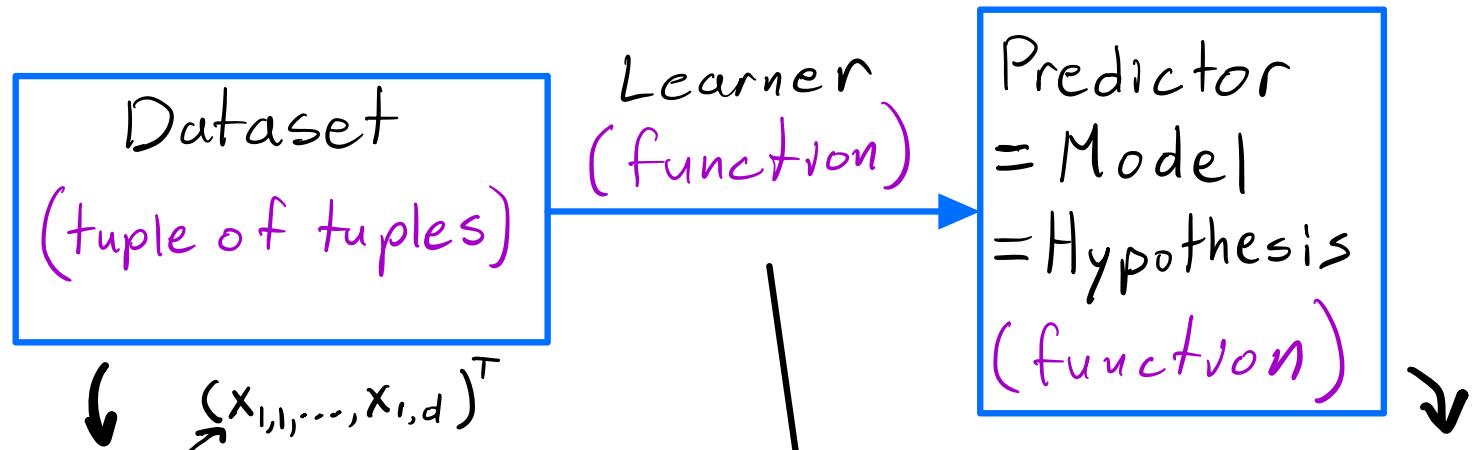
$X \sim P \rightarrow$ "The r.v. X is sampled according
to the distribution P "

typo : RV section $Y \in \mathbb{Y} \in \{0,1\}$

- two ways to describe a discrete random experiment
- cont. r.v. statements
- question in past lecture : infinite outcome set possible ?

Motivation

Supervised Learning: Learning from a randomly sampled batch of labeled data



$D = ((\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n))$
 $D \in (\mathcal{X} \times \mathcal{Y})^n$
 $(\vec{x}_i, y_i) \sim P_{\vec{x}, y}$
independent for all $i \in \{1, \dots, n\}$

D n feature-label pairs
 \mathcal{X} set of features
 \mathcal{Y} set of labels/targets
Ex: $\mathcal{X} = \mathbb{R}^2$, $\mathcal{Y} = \mathbb{R}$

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

f a function from features to labels

$$\underline{\text{Ex}} \quad f(x) = 2x + 1, \mathcal{X} = \mathbb{R}$$

$$\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\}$$

\mathcal{A} a function from datasets to predictors

$$\underline{\text{Ex:}} \quad \mathcal{A}(D) = f \text{ where } f: \mathcal{X} \rightarrow \mathcal{Y}$$

$$f(x) = \begin{cases} y_i & \text{if } x = x_i \text{ for some } i \in \{1, \dots, n\} \\ 0 & \text{otherwise} \end{cases}$$

Probability

Note: Humans have a bad intuition when it comes to randomness
-Thinking Fast and Slow
by: Daniel Kahneman

- Random Variables
- Calculating probabilities using pmf and pdf
- Multivariate random variables
 - Conditional and marginal probabilities
- Representing random features, labels, and datasets
- Functions of random variables
- Expectation and variance

Warning: If some things seem informal, it is likely because we would need tools from Measure Theory, which we will not cover in this course.

Experiment: A process that generates an uncertain outcome

Ex: flipping a coin, rolling a dice

Outcome Space/Set: The set of all outcomes from the experiment

Ex: $y = \{0, 1\}$ Heads
Tails flipping a coin

$\chi = \{1, 2, 3, 4, 5, 6\}$ rolling a dice

$[0, 900]$ amount of a chemical in a wine
 \mathbb{R} measurement error

The outcome space / set can be either a

- 1) countable set or
- 2) uncountable set

→ cardinality finite or
countably infinite

→ cardinality is uncountably infinite

Event: A subset of the outcome space (imprecise)

Ex.: Outcome space: $Y = \{0, 1\}$

Events: $\{0\}, \{1\}, \{0, 1\} = Y, \emptyset$

an outcome is a single element of the outcome set

Ex.: Outcome space: $X = \{1, 2, 3, 4, 5, 6\}$

Events: $\emptyset, X, \{1\}, \{1, 3, 5\}, \{2, 4, 6\}, \{1, 2, 3\}, \dots$

Ex.: Outcome space: $[0, 900]$

Events: $\emptyset, [0, 900], [0, 4], [1, 2] \cup [7, 30], \dots$

Probability Distribution: A function P defining the likelihood of each event (and satisfying certain properties)

P : event space / set $\rightarrow [0, 1]$

a set containing all the events

A complicated set (σ -algebra) that we will not define

Properties: (imprecise)

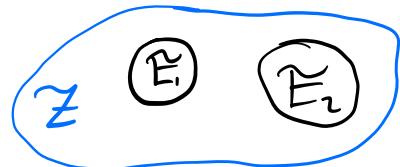
Outcome space: \mathcal{Z}

1. $P(\mathcal{Z}) = 1$

2. If $E_1 \subset \mathcal{Z}, E_2 \subset \mathcal{Z}$ and $E_1 \cap E_2 = \emptyset$, then

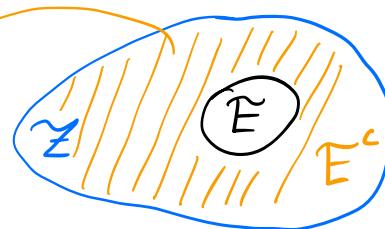
$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

disjoint
↓



Ex: (of property 2.)

Events: E, E^c



$$E \cap E^c = \emptyset, E \cup E^c = \mathcal{Z}$$

$$\underbrace{P(E \cup E^c)}_{= P(\mathcal{Z})} = P(E) + P(E^c)$$

$$= P(\mathcal{Z}) = 1$$

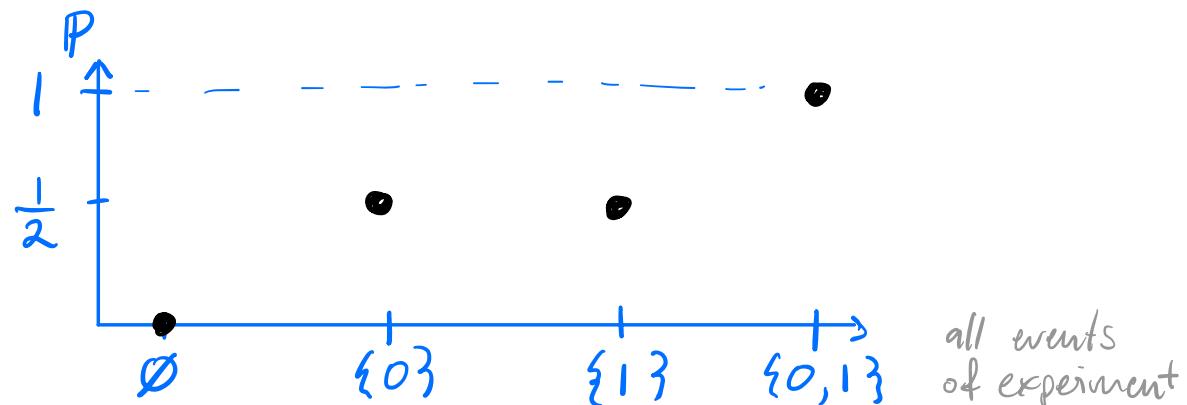
rearranging:

$$P(E) = 1 - P(E^c)$$

Ex: Outcome space: $\mathcal{Y} = \{0, 1\}$

$$P(\emptyset) = 0, P(\mathcal{Y}) = 1, P(\{0\}) = \frac{1}{2}, P(\{1\}) = \frac{1}{2}$$

each event of an experiment is associated with a probability

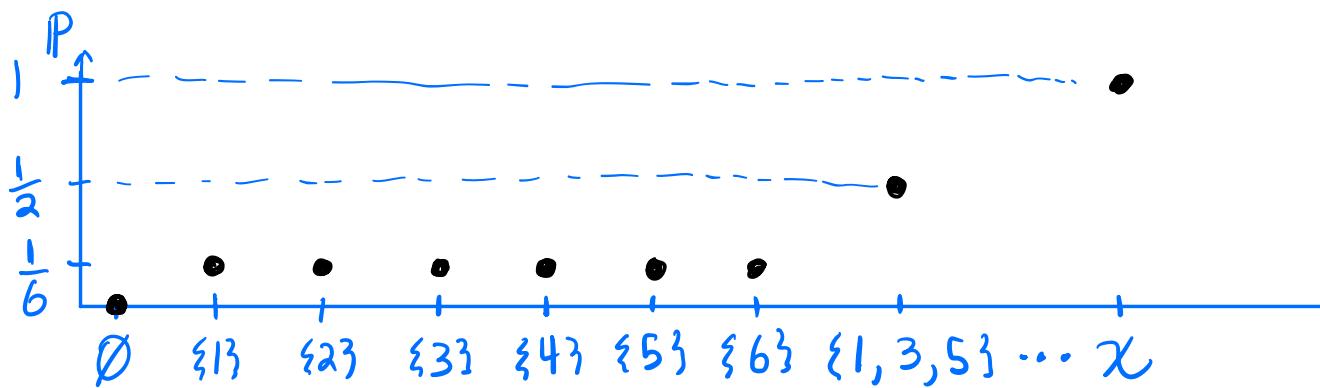


Ex: Outcome space: $\chi = \{1, 2, 3, 4, 5, 6\}$

$$P(\{1\}) = P(\{2\}) = \dots = P(\{6\}) = \frac{1}{6}$$

$$P(\{1, 3, 5\}) = P(\{1\}) + P(\{3\}) + P(\{5\}) = \frac{1}{2}$$

$$P(\chi) = 1$$



Random Variables

Random Variable (r.v.): (imprecise) A variable that takes a value based on the outcome of an experiment, and is associated with a probability distribution. The value can be any random outcome.

Ex: $X \in \chi = \{1, 2, 3, 4, 5, 6\}$ with P from prev. example
 $Y \in \mathcal{Y} = \{0, 1\}$ with P from prev. example
 $Z \in \{H, T\}$

A random variable is actually a function (satisfying certain properties) from one outcome space to another outcome space. Ex $X(T) = 0, X(H) = 1$.

It will not be necessary to know this for this course

Probability Distributions with r.v.

Ex: Outcome space: \mathcal{X} r.v.: $X \in \mathcal{X}$

$$P(\{1, 3, 5\}) \stackrel{\text{def}}{=} P(X \in \{1, 3, 5\})$$

"the r.v. X is an element
of the event $\{1, 3, 5\}$ "

$$P(\{4, 5, 6\}) = P(X \in \{4, 5, 6\}) = P(X \geq 4)$$

$$P(\{4\}) = P(X \in \{4\}) = P(X = 4)$$

Notation: $Z \sim P$ "Z is sampled according to distribution P"

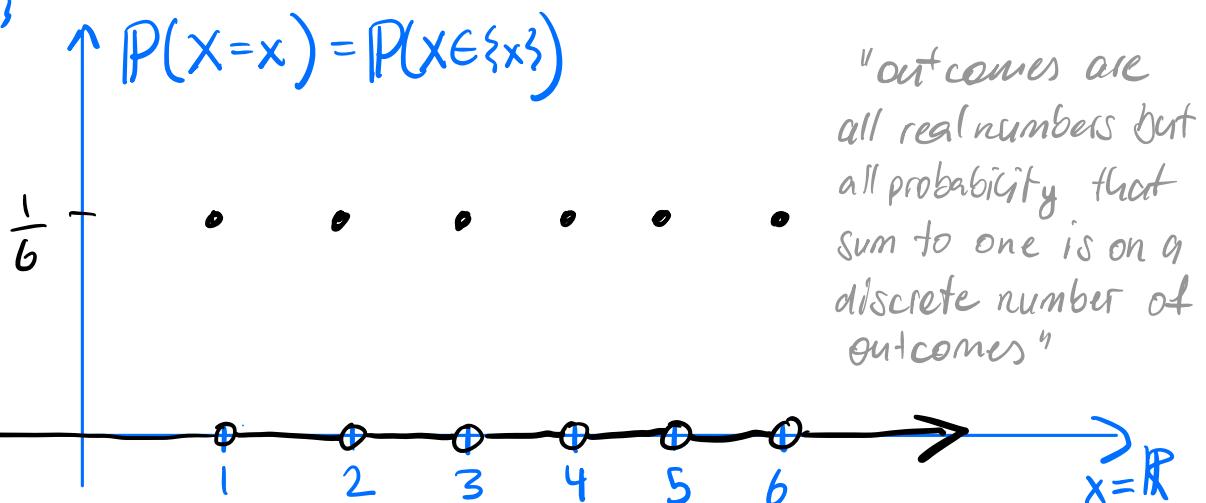
Discrete r.v.: A r.v. that takes values from:

- A countable outcome space, or
- an uncountable outcome space, but there is a countable event that has probability 1

Ex: $Y \in \mathcal{Y} = \{0, 1\}$, $X \in \mathcal{X} = \{1, 2, 3, 4, 5, 6\}$, $Z \in \mathbb{N}$ ← two ways to describe same discrete experiment

Ex: $X \in \mathbb{R}$ where $P(X=1) = \dots = P(X=6) = \frac{1}{6}$ ←

Probability 1 so $P(X \in \{1, 2, 3, 4, 5, 6\}) = 1$
 for countable and $P(\mathbb{R} \setminus \{1, 2, 3, 4, 5, 6\}) = 0$
 event $\{1, 2, 3, 4, 5, 6\}$



Note: You can always take a r.v. defined on a countable outcome space and define it on a larger uncountable outcome space by setting the probability of the event containing all the new outcomes to zero

Continuous r.v.: A r.v. that takes values from:

- an uncountable outcome space and the probability of any single outcome is zero

Ex: $Z \in [0, 900]$ and $P(Z=z) = P(Z \in \{z\}) = 0$ for all $z \in [0, 900]$
but $P(Z \in [0, 900]) = 1$ =

Ex: $Z \in \mathbb{R}$ and $P(Z=z) = P(Z \in \{z\}) = 0$ for all $z \in \mathbb{R}$
but $P(Z \in \mathbb{R}) = 1$ =

Calculating Probabilities

Motivation: It is hard to define the values of a probability distribution P for all the events

Probability Mass Function (pmf): A function $p: \mathcal{Z} \rightarrow [0, 1]$

where \mathcal{Z} is a discrete outcome space and $\sum_{z \in \mathcal{Z}} p(z) = 1$.

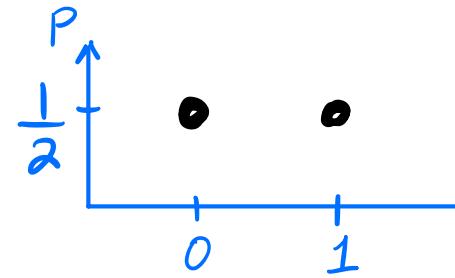
The probability of an event $E \subset \mathcal{Z}$ is:

$$P(Z \in E) \stackrel{\text{def}}{=} \sum_{z \in E} p(z)$$

where $Z \in \mathcal{Z}$

Ex: Outcome space: \mathcal{Y}

$$P(0) = \frac{1}{2}, P(1) = \frac{1}{2}$$



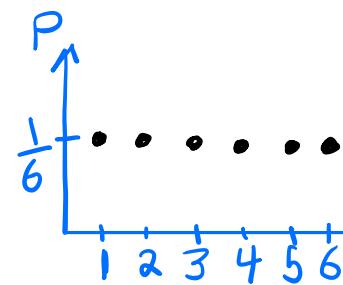
$$P(Y \in \{0, 1\}) = \sum_{y \in \{0, 1\}} P(y) = P(0) + P(1) = 1$$

$$P(Y=0) = P(Y=1) = P(Y \in \{1\}) = \sum_{y \in \{1\}} P(y) = P(1) = \frac{1}{2}$$

$$P(Y \in \emptyset) = \sum_{y \in \emptyset} P(y) = 0$$

Ex: Outcome space: \mathcal{X}

$$P(1) = P(2) = \dots = P(6) = \frac{1}{6}$$



$$P(X \in \{1, 3, 5\}) = \sum_{x \in \{1, 3, 5\}} P(x) = P(1) + P(3) + P(5) = \frac{1}{2}$$

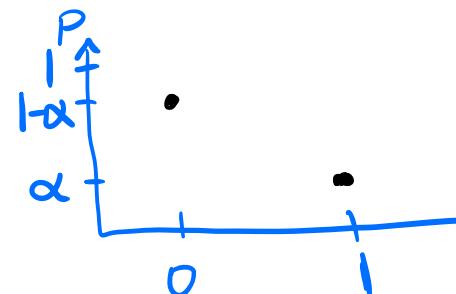
Discrete Probability Distributions with special names:

Bernoulli distribution (parameter: $\alpha \in [0, 1]$):

Outcome space: $\{0, 1\}$

$$\text{pmf: } P(1) = \alpha, P(0) = 1 - \alpha$$

Distribution $P = \text{Bernoulli}(\alpha)$



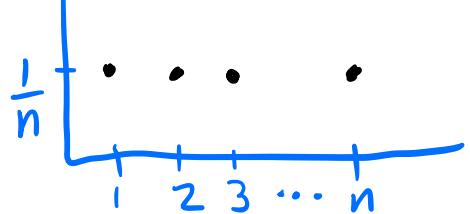
$$P(Z=1) = P(Z \in \{1\}) = P(1) = \alpha \quad Z \in \{0, 1\} \text{ is a "Bernoulli r.v."}$$

Discrete Uniform Distribution (parameter: n):

Outcome space: $\{1, 2, \dots, n\}$

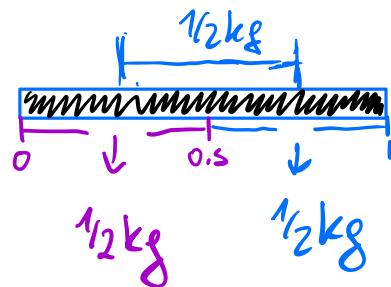
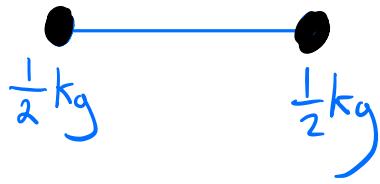
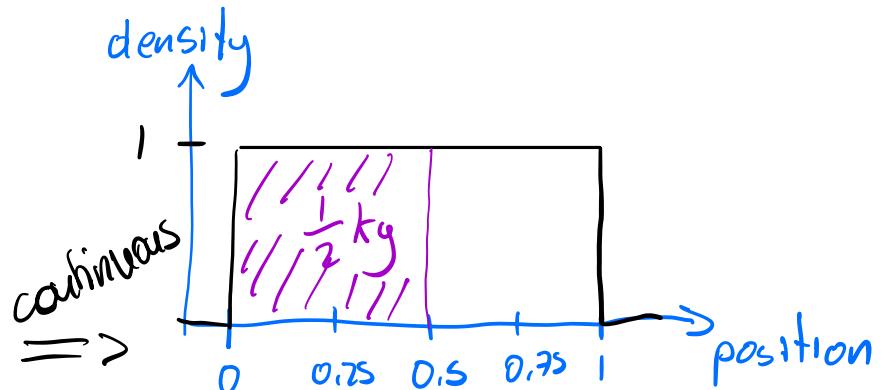
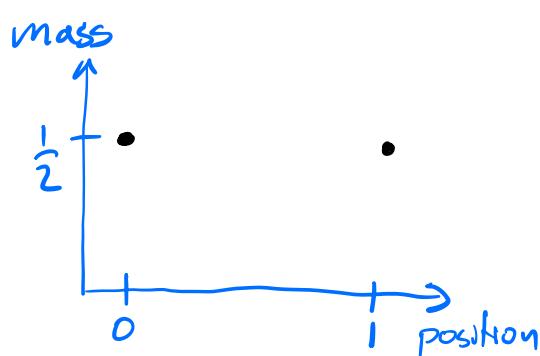
P_{\uparrow}

pmf: $p(1) = p(2) = \dots = p(n) = \frac{1}{n}$



Distribution $P = \text{Uniform}(n)$

Intuition with a rod in physics



Probability Density Function (pdf): a function $p: \mathbb{Z} \rightarrow [0, \infty]$

where \mathbb{Z} is an uncountable outcome space and $\int_{\mathbb{Z}} p(z) dz = 1$

The probability of an event $E \subset \mathbb{Z}$ is:

$$P(Z \in E) \stackrel{\text{def}}{=} \int_E p(z) dz$$

$$P(Z = z) = P(Z \in \{z\}) = 0$$

where $z \in \mathbb{Z}$

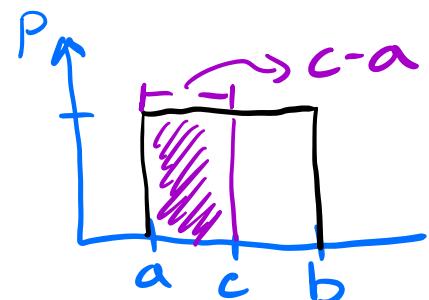
Continuous Probability Distributions with Special names:

Continuous Uniform Distribution (parameters: $a \in \mathbb{R}, b \in \mathbb{R}$):

Outcome space: $[a, b]$

pdf: $p(z) =$

Distribution $\mathbb{P} = \text{Uniform}(a, b)$



$$\mathbb{P}(a \leq z \leq c) = \mathbb{P}(z \in [a, c]) = \int_a^c p(z) dz =$$

where $a \leq c \leq b$

Gaussian/Normal Distribution (parameters: $\mu \in \mathbb{R}, \sigma^2 > 0$):

Outcome space: \mathbb{R}

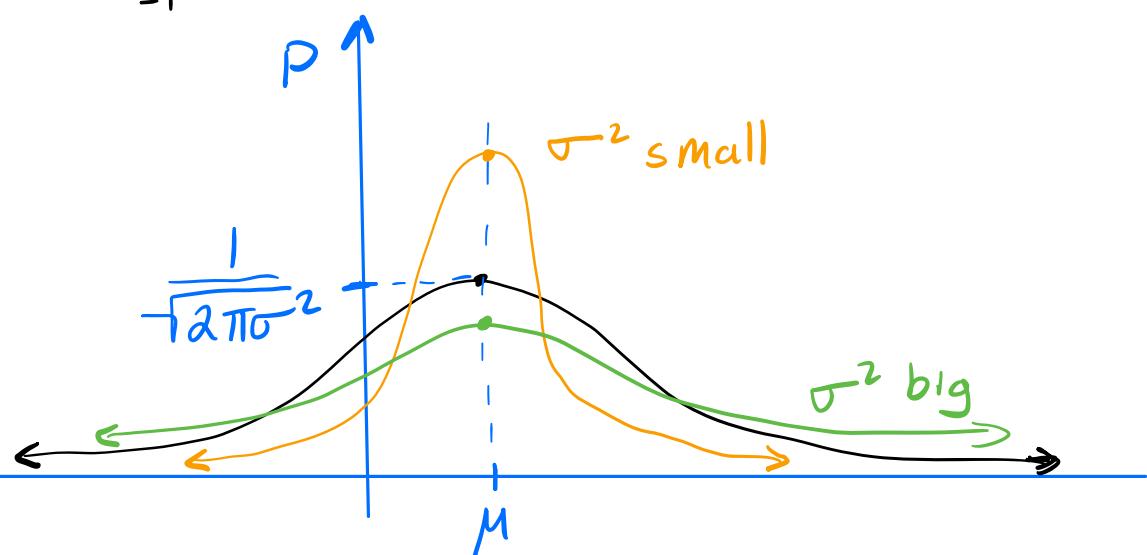
$$\text{pdf: } p(z) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(z-\mu)^2\right)$$

standard normal dist.
 $\rightarrow N=0, \sigma^2=1$

Distribution $\mathbb{P} = N(\mu, \sigma^2) = \text{Gaussian}(\mu, \sigma^2)$

$\sigma \rightarrow$ standard deviation

$$\underline{\text{Ex}} \quad \mathbb{P}(-1 \leq z \leq 1) = \int_{-1}^1 p(z) dz$$

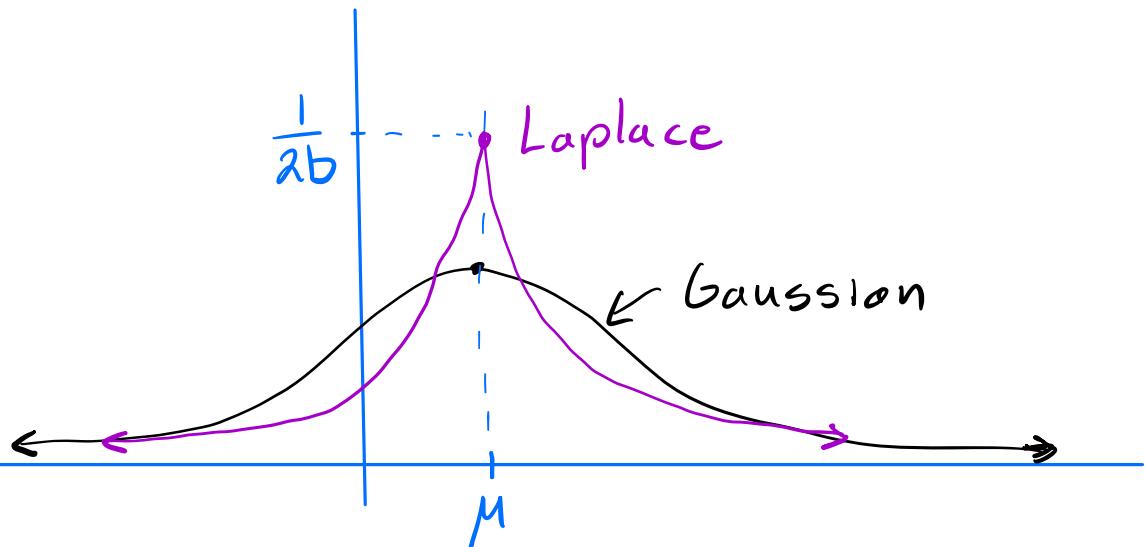


Laplace Distribution (parameters: $\mu \in \mathbb{R}$, $b > 0$):

Outcome space: \mathbb{R}

pdf: $p(z) = \frac{1}{2b} \exp\left(-\frac{1}{b}|z-\mu|\right)$

Distribution $P = \text{Laplace}(\mu, b)$



Multivariate Random Variables

Motivation: To be able to talk about the probability of different types of events (outcomes of different experiments) at the same time!

- Ex: The probability of getting heads and rolling a 3
- The probability of a wine containing 2.5mg of one chemical and 4mg of another chemical
- The probability of a house having 4 rooms and 2 washrooms and being less than 10min from a university
- The probability of being young and having arthritis

Multivariate Random Variable: A tuple of more than one random variable

Ex: (Flipping 2 coins) Heads

Outcome space: $\mathcal{X} = \{0, 1\} \times \{0, 1\} = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$

r.v.: $X = (X_1, X_2) \in \mathcal{X}$

(Collecting the info of one house (ex: # of rooms, age))

Outcome space: $\mathcal{X} = \mathbb{N} \times [0, \infty)$ age

r.v.: $X = (X_1, X_2) \in \mathcal{X}$ # of rooms

(Collecting the info of one house and its price)

Outcome space: $\mathcal{Z} = (\mathbb{N} \times [0, \infty)) \times [0, \infty)$ age

r.v.: $Z = (X, Y)$ # of rooms Price

$= ((X_1, X_2), Y)$

Calculating Joint Probabilities

Ex: (If you have arthritis and if you are young or old)

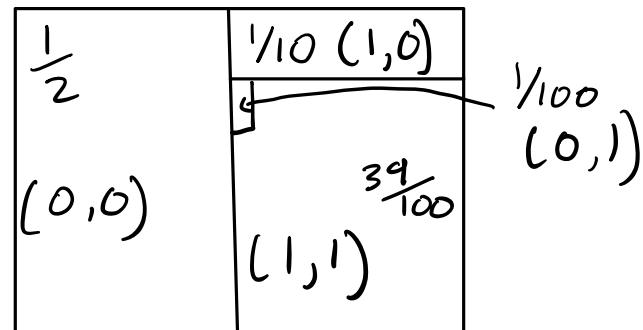
Outcome space: $\mathcal{Z} = \{0, 1\} \times \{0, 1\}$ Arthritis
Old ↓ ↓ Young No arthritis
Young ↑ Arthritis

r.v.: $Z = (X, Y) \in \mathcal{Z}$

pmf: $p: \mathbb{Z} \rightarrow [0, 1]$

Not based on real data $\left\{ \begin{array}{l} p((0,0)) = p(0,0) = \frac{1}{2}, \quad p(0,1) = \frac{1}{100} \\ p(1,0) = \frac{1}{10}, \quad p(1,1) = \frac{39}{100} \end{array} \right.$

	Y	
X	0	1
0	$\frac{1}{2}$	$\frac{1}{100}$
1	$\frac{1}{10}$	$\frac{39}{100}$



$$\sum_{z \in \mathbb{Z}} p(z) = \sum_{x \in X} \sum_{y \in Y} p(x,y) = \frac{1}{2} + \frac{1}{100} + \frac{1}{10} + \frac{39}{100} = 1 \quad (\text{for a valid pmf})$$

What is the probability of being young (i.e. $X=0$)?

$$\mathbb{E} = \{(0,0), (0,1)\} = \{0\} \times \{0, 1\} \subset \mathbb{Z}$$

$$\begin{aligned} P(X=0, Y \in \{0, 1\}) &= P(Z \in \mathbb{E}) = \sum_{z \in \mathbb{E}} p(z) \\ &= \sum_{x \in \{0\}} \sum_{y \in \{0, 1\}} p(x, y) \\ &= p(0,0) + p(0,1) \\ &= \frac{1}{2} + \frac{1}{100} = \frac{51}{100} \end{aligned}$$

Marginal Distribution: The distribution over a subset of random variables

Ex: Continuing with the arthritis example

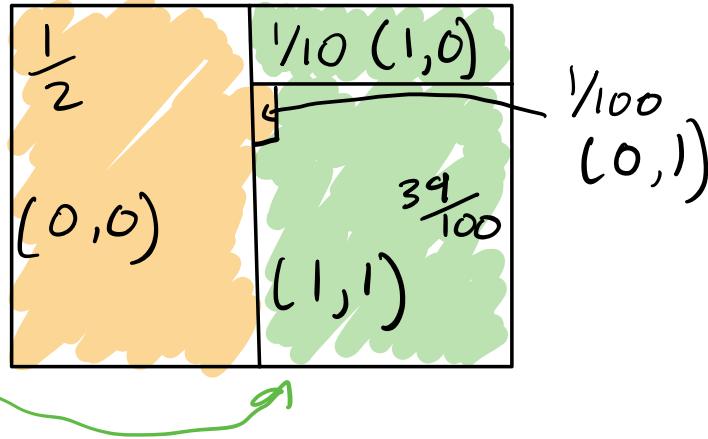
Marginal Distribution: $P_x(X \in E_x)$ where $E_x \subset \mathcal{X}$

Marginal pmf: $P_x: \mathcal{X} \rightarrow [0, 1]$, $P_x(x) = \sum_{y \in Y} P(x, y)$

$$P_x(X=0) = P_x(0) = \sum_{x \in \{0\}} \sum_{y \in Y} P(0, y) = \frac{51}{100}$$

$$P_x(X=1)$$

$$= \frac{1}{10} + \frac{39}{100} = \frac{49}{100}$$



Discrete r.v. X_1, \dots, X_d

$X = (X_1, \dots, X_d) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_d = \mathcal{X}$, $P: \mathcal{X} \rightarrow [0, 1]$

$P_{X_i}: \mathcal{X}_i \rightarrow [0, 1]$, $i \in \{1, \dots, d\}$

"Sum over every element except x_i "

Marginal pmf:

$$P_{X_i}(x_i) \stackrel{\text{def}}{=} \sum_{x_1 \in \mathcal{X}_1} \dots \sum_{x_{i-1} \in \mathcal{X}_{i-1}} \sum_{x_{i+1} \in \mathcal{X}_{i+1}} \dots \sum_{x_d \in \mathcal{X}_d} P(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d)$$

$$P_{X_i}(X_i \in E_i) = \sum_{x_i \in E_i} P_{X_i}(x_i) \quad \text{where } E_i \subset \mathcal{X}_i$$

Continuous r.v. X_1, \dots, X_d

$X = (X_1, \dots, X_d) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_d = \mathcal{X}$, $P: \mathcal{X} \rightarrow [0, \infty)$

$P_{X_i}: X_i \rightarrow [0, \infty)$, $i \in \{1, \dots, d\}$

Marginal
pdf:

$$P_{X_i}(x_i) = \int_{x_1} \int_{x_{i-1}} \int_{x_{i+1}} \int_{x_d} P(x_1, \dots, x_d) dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_d$$

Distribution: $P_{X_i}(X_i \in E_i) = \int_{E_i} P_{X_i}(x_i) dx_i$ where $E_i \subset X_i$

Conditional Distribution: Probability of
a r.v. given info about another r.v.

Ex: Probability that I have arthritis given I am young

Let r.v. $= Y \in \mathcal{Y}, X \in \mathcal{X}$

Discrete Y for any $x \in \mathcal{X}$ that $P_X(x) \neq 0$

$P_{Y|X=x}: Y \rightarrow [0, 1]$, $P_{Y|X=x}(y) = P_{Y|X}(y|x)$

conditional
pmf: $P_{Y|X}(y|x) \stackrel{\text{def}}{=} \frac{P(y, x)}{P_X(x)}$ implies $\sum_{y \in \mathcal{Y}} P_{Y|X}(y|x) = 1$

Distribution: $P_{Y|X}(Y \in E_Y | X=x) \stackrel{\text{def}}{=} \sum_{y \in E_Y} P_{Y|X}(y|x)$ where $E_Y \subset \mathcal{Y}$

Continuous Y for any $x \in \mathcal{X}$ that $P_X(x) \neq 0$

$P_{Y|X=x}: Y \rightarrow [0, \infty)$, $P_{Y|X=x}(y) = P_{Y|X}(y|x)$

conditional
pdf:

$$p_{Y|X}(y|x) \stackrel{\text{def}}{=} \frac{P(y,x)}{P_X(x)}$$

implies $\int_y p_{Y|X}(y|x) dy = 1$

Distribution:

$$P_{Y|X}(Y \in E | X=x) \stackrel{\text{def}}{=} \int_E p_{Y|X}(y|x) dy \quad \text{where } E \subset Y$$

Product Rule:

$$p(x,y) = p(x|y)p(y) = p(y|x)p(x)$$

More generally:

$$p(x_1, x_2, \dots, x_d) = p(x_d | x_1, \dots, x_{d-1}) \dots p(x_3 | x_1, x_2) p(x_2 | x_1) p(x_1)$$

Bayes' Rule:

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Note: Sometimes the subscripts are not used for marginal and conditional distributions when it is clear from the context

$$p(x,y) = p_{x,y}(x,y)$$

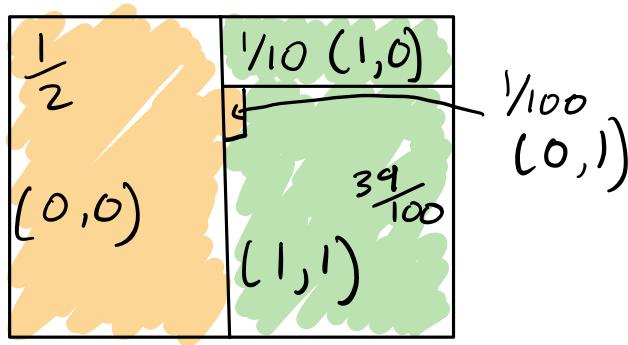
$$p(x) = P_X(x)$$

$$p(y|x) = P_{Y|X}(y|x)$$

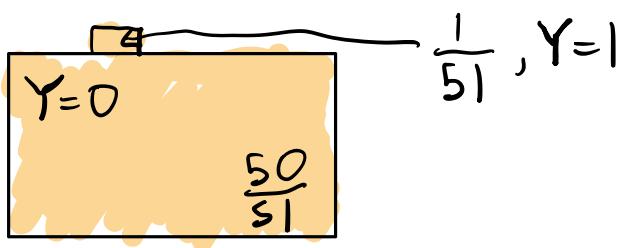
Ex: Probability that I have arthritis given I am young

$$P(Y=1 | X=0) = \sum_{Y \in \{1\}} P_{Y|X}(y|0)$$

Arthritis Young = $P_{Y|X}(1|0)$



↓ condition on
being young



$$= \frac{P_{X|Y}(0,1)}{P_X(0)}$$

$$= \frac{P(0,1)}{P(0,1) + P(0,0)}$$

$$= \frac{\frac{1}{100}}{\frac{1}{100} + \frac{1}{2}}$$

$$= \frac{\frac{1}{100}}{\frac{51}{100}} = \frac{1}{51}$$

Ex: Probability of being young given I have arthritis

$$P(X=0|Y=1) = P_{X|Y}(0|1)$$

Bayes' Rule $\Rightarrow \frac{P_{Y|X}(1|0) P_X(0)}{P_Y(1)}$

$$= \frac{\frac{1}{51} \frac{51}{100}}{\frac{40}{100}} = \frac{1}{40}$$

Independence: Changing the value of one r.v. doesn't affect the probability of another r.v.

r.v. X, Y are independent if: $P(x, y) = p_x(x)p_y(y)$

Since $P(x, y) = P(x|y)p(y) = P(x)p(y|x) = P(x)p(y)$

independence implies: $P(x|y) = P(x)$, $P(y|x) = P(y)$

More generally:

X_1, X_2, \dots, X_d are independent if: $P(x_1, \dots, x_d) = P(x_1) \cdots P(x_d)$

Similarly for distributions:

r.v. X, Y are independent if: $P(X \in E_x, Y \in E_y) = P(X \in E_x)P(Y \in E_y)$

Ex: X, Y are not independent for Arthritis ex

$$P(0, 1) = \frac{1}{100} \neq P_x(0)P_y(1) = \frac{51}{100} \cdot \frac{40}{100} = 0.204$$

Ex: $X_1, X_2 \in \{0, 1\}$ are flips of two different fair coins

$$P(x_1, x_2) = \frac{1}{4} \text{ for all } x_1 \in \mathcal{X}_1, x_2 \in \mathcal{X}_2$$

$$P_{x_1}(x_1)P_{x_2}(x_2) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

		x_2	
		H	T
x_1	H	$\frac{1}{4}$	$\frac{1}{4}$
	T	$\frac{1}{4}$	$\frac{1}{4}$